



(19)

Europäisches Patentamt
European Patent Office
Office européen des brevets



(11)

EP 0 717 355 B1

(12)

EUROPEAN PATENT SPECIFICATION

(45) Date of publication and mention
of the grant of the patent:
26.09.2001 Bulletin 2001/39

(51) Int Cl.⁷: **G06F 9/46, G06F 9/44**

(21) Application number: **95308908.3**

(22) Date of filing: **07.12.1995**

(54) Parallel processing system and method

Paralleles Verarbeitungssystem und -verfahren

Système et méthode de traitement parallèle

(84) Designated Contracting States:
DE FR GB

(30) Priority: **12.12.1994 US 353590**

(43) Date of publication of application:
19.06.1996 Bulletin 1996/25

(73) Proprietor: **NCR International, Inc.**
Dayton, Ohio 45479 (US)

(72) Inventors:

- **Kandasamy, David R.**
San Ramon, CA 94583 (US)
- **Heying, Douglas W.**
Richmond, Surrey TW10 6HQ (GB)
- **Catozzi, John R.**
Redondo Beach, CA 90278 (US)

(74) Representative: **Cleary, Fidelma et al**
International IP Department
NCR Limited
206 Marylebone Road
London NW1 6LY (GB)

(56) References cited:
EP-A- 0 602 773

- **SOFTWARE PRACTICE & EXPERIENCE**, vol. 21,
no. 10, 1 October 1991, pages 989-1013,
XP000297890 AUSTIN P ET AL: "THE DESIGN OF
AN OPERATING SYSTEM FOR A SCALABLE
PARALLEL COMPUTING ENGINE"
- **FUTURE GENERATIONS COMPUTER**
SYSTEMS, vol. 8, no. 1 / 03, 1 July 1992, pages
93-109, **XP000343296 STEINER P: "EXTENDING**
MULTIPROGRAMMING TO A DMPP"
- **PATENT ABSTRACTS OF JAPAN** vol. 014, no.
345 (P-1083), 26 July 1990 & **JP-A-02 123455**
(**HITACHI LTD**), 10 May 1990,

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

Description

[0001] This invention relates to parallel processing systems.

[0002] Parallel processing systems are frequently comprised of an operating system and arrays of individual computers (i.e., processor nodes), each with their own central processing unit (CPU), memory, and data storage unit.

Tasks are executed in parallel by utilizing each processor node.

[0003] During the execution of a task, a body of work is divided into multiple threads. A thread is a stream of instructions executed by the computer on behalf of a task. As an analogy, a task such as an orchestra performing a symphony can be decomposed into many threads which would be the individual musicians, each playing their part. Typically, in a parallel computer system, each thread is allocated to a different processor node. Each of these threads is then executed in parallel at their respective separate nodes. For instance, three threads can occupy and execute simultaneously on three different nodes at the same time. One of the advantages of this parallel processing technique is accelerated execution times.

[0004] Although the conventional form of parallel processing has merits, there are shortcomings. Conventional parallel processing techniques may result in an inefficient use of the available hardware. For example, if a processor node has multiple storage devices (i.e., disk drives) attached, then a single thread of execution might access only one of those devices at a time, leaving the other storage devices underutilized or even idle. In addition, if a parallel system is constructed out of multiprocessor nodes, then a single thread per node might not utilize all of the available CPUs in the node.

[0005] SOFTWARE PRACTICE AND EXPERIENCE, Vol. 21, No.10, 1 October 1991 pages 989-1013, XPOO 297890 Austin P et al: "Design of an operating system for a scalable parallel processing engine" discloses a computer system including a plurality of processor nodes and an interconnect network coupled to the processor nodes for transmitting messages between processor nodes. The document further discloses that each of the processor nodes comprises a plurality of virtual processor means for performing a plurality of threads simultaneously, whereby each virtual processor means performs one of the threads. In this device each of the processor nodes comprises operating system means for the assigning of virtual processing means to processor nodes at system start up.

[0006] Patent Abstracts of Japan volume 14 no. 345 (P-1083) 26 July 1990 and JP-A-02123455 (Hitachi Ltd), 10 May 1990 disclose a process for the dynamic reconstitution of a system by separating a failed processor element from a system and allocating the function of the elements to another normal processor element as a virtual processor element. A translation table identify-

ing locations of virtual processor means and processor nodes on an interconnect network is disclosed in order to enable the reconstitution of the system.

[0007] It is an object of the present invention to provide a more efficient usage of parallel computer systems.

[0008] Therefore, according to one aspect of the present invention, there is provided a computer system, including a plurality of processor nodes, each of the processor nodes having a plurality of data storage devices connected thereto; and an interconnect network, coupled to the processor nodes for transmitting messages between processor nodes, each of the processor nodes comprising a plurality of virtual processor means for performing a plurality of threads simultaneously, wherein each virtual processor means performs one of the threads, the threads being streams of instructions executed on behalf of a task, and each of the processor nodes comprising operating system means for assigning the virtual processor means to processor nodes at system startup, characterized in that each of the threads accesses data on a different data storage device simultaneously with the other threads; and further comprising means for creating and storing a translation table identifying locations of virtual processor means and processor nodes on the interconnect network at system startup.

[0009] According to another aspect of the present invention, there is provided a method of executing instructions in a parallel computer system comprised of a plurality of processor nodes, each of the processor nodes having a plurality of data storage devices connected thereto, and an interconnect network, coupled to the processor nodes, for transmitting messages between processor nodes, including the steps of: utilizing a plurality of virtual processor means to perform multiple threads simultaneously, wherein each virtual processor means comprises an operating environment in one of the processor nodes for performing one of the threads, the threads being streams of instructions executed on behalf of a task, and assigning each of the virtual processor means to one of the processor nodes at system startup, characterized in that each of the threads accesses data on a different data storage device simultaneously with the other threads; and further comprising creating and storing a translation table identifying locations of virtual processor means and processor nodes on the interconnect network at system startup.

[0010] One embodiment of the present invention will now be described by way of example, with reference to the accompanying drawings, in which:-

Figure 1 is a diagram of a typical parallel processing system;

Figure 2 is a representation of disk storage usage without vprocs;

Figure 3 is a representation of disk storage usage with vprocs;

Figures 4A and 4B illustrate how vprocs communicate with each other using addresses that are vproc-specific;

Figure 5 is a representation of a clique, a collection of processor nodes connected to the same data storage; and

Figure 6 illustrates that the vproc configuration can be held constant, even though one or more processor nodes are down.

[0011] Figure 1 shows a typical parallel processing system. An interconnect network 10 is the means by which the processor nodes 12 are grouped together and communicate with each other. Each node may have multiple storage devices 14. A processor node 12 may also have multiple CPUs used for tightly coupled or symmetric multiprocessing. In a typical parallel processing system, a task is divided into multiple threads of execution (or elements of work), usually with one thread per processor node 12. Each of these threads is then executed in parallel (concurrently) by the processor nodes 12. However, if a node 12 is constructed out of multiple CPUs, a single thread per node might not occupy all the CPUs.

[0012] Figure 2 is a representation of a single node 12 without vprocs and with multiple storage devices attached 14. This figure illustrates that a single thread 16 of execution might access three devices 14 sequentially, but only one at a time, leaving the other available storage devices 14 idle.

[0013] Figure 3 is a representation of a single node 12 with multiple threads 16 and multiple storage units 14. Each thread 16 is encapsulated in an operating environment termed a Virtual Processor (vproc) 18. Each vproc 18 is given its own private logical disk space, called a Virtual Disk (vdisk), which may actually be one or more data storage devices 14. Using the vproc concept, the degree of parallelism provided by the parallel system may be increased to include multiple concurrently executing threads 16 per node 12, so that there is one thread 16 per storage device 14, rather than one thread 16 per node 12 as shown in Figure 2. The combination of vprocs 18 and vdisks allows the threads 16 to execute simultaneously and not just alternate time slices of the node 12. Similarly, if a processor node 12 has multiple CPUs, each individual thread 16 may utilize a different CPU for enhanced parallelism. The vproc 18 also provides a level of isolation from other threads 16 operating in other vprocs 18 on the same node 12. Moreover, the vprocs 18 provide a degree of location transparency, in that vprocs 18 communicate with each other using addresses that are vproc-specific, rather than processor-node-specific, i.e., a vproc 18 can be addressed without needing to know what node it's executing on and vprocs 18 on different nodes 12 talk to each other using the same mechanisms as vprocs on the same node 12. Further, vprocs 18 facilitate redundancy by providing a level of isolation/abstraction between the

physical node 12 and the thread 16 of execution, thereby allowing all threads 16 to run even when some nodes 12 have failed or are shut down. The result is increased system utilization and fault tolerance.

5 **[0014]** Figure 4A illustrates the level of system software at which the preferred embodiment of the vproc concept can be implemented. Traditional partitioning of computer systems has been done at a very low level, where address spaces are divided, and interrupts and
10 other low-level CPU operations are partitioned and isolated. The present invention implements vprocs 18 at the operating system data structure level and in its interconnect 10 software. By providing "virtualism" at the data structure level, the partitioning and isolation of
15 vprocs 18 can be provided at a much lower performance cost than traditional virtual machine implementations.

[0015] The operating system (OS) image in each node 12 keeps track of which vproc 18 each task belongs to. Whenever that task makes an operating system call to perform some function, the OS keeps track of addresses/tags/handles for the entity on a vproc 18 basis. In other words, the task will get back a mailbox with an address that is unique within that vproc 18. Other vprocs 18 operating on the same node 12 may also have
20 that address, but since the OS keeps track of which tasks belong to which vproc 18, all those mailboxes are logically distinct.

[0016] The OS also assigns vprocs 18 to processor nodes 12 at system startup time, preferably based on a
30 load-balancing algorithm, although other allocation algorithms may be used as well. In addition, the OS can re-assign vprocs 18 to processor nodes 12 whenever additional processor nodes 18 become available, i.e., after they have been repaired and/or rebooted. At system startup time, a translation table of vprocs 18 and the
35 nodes 12 on which they reside is created and stored in an interconnect driver 10. The translation table of vprocs 18 and the nodes 12 on which they reside is updated whenever vprocs 18 are moved between nodes 12 for
40 whatever reason.

[0017] When a message is sent to a mailbox in a particular vproc 18, the interconnect driver 10 looks up the vproc/node translation in the table, sends the message to the designated node 12, and the recipient node 12
45 then routes it locally to the appropriate mailbox in the proper vproc 18. Thus, vprocs 18 communicate with each other using addresses that are vproc-specific rather than node-specific. Specifically, a vproc 18 can be addressed without knowing what processor node 12 it is executing on.

[0018] For example, referring to Figure 4B, a task in vproc0 22 may decide to send a message to mbx1 in mailbox 24 in vproc3 26. The message is prepared and sent to the interconnect driver 28 on the local processor
55 node 30. The driver 28 uses the destination vproc3 26 to look up the hosting processor node in its translation table 34, in the example thus identifying node1 32. The message is sent across the interconnect network 36 to

node 32. The interconnect network driver 28 in node 32 receives the message and routes it to mbx1 in mailbox 24 in vproc3 26. As a result of the vproc-specific addresses, a task located anywhere in the computer system may send a message to any other task located anywhere in the system. Therefore, referring to Figure 4A, vprocs 18 on the same node 12 communicate with each other using the same mechanism as vprocs 18 on different nodes 12.

[0019] Figure 5 illustrates the concept of a clique. A clique is a collection of processor nodes 12 connected to one or more data storage devices 14. Cliques provide a clustering of processor nodes 12 for increased availability with the side benefit of common shared access to data storage devices 14. Vdisks can be accessed from any processor node 12 within a clique, thereby allowing a vproc 18 to run on any node 12 within the clique. By grouping processor nodes 12 together and attaching them to a common disk storage devices 14, it can be shown that the likelihood of at least one of the processor nodes 12 in the group being functional is orders of magnitude higher than the likelihood of all of the processor nodes in a group being functional. For example, in a three node 12 clique made up of nodes 12 with a 30,000 hr. Mean Time Between Failure (MTBF), the average time between failure of nodes 12 within a clique is 10,000 hr. (30,000 hr./3 nodes per clique). However, given a Mean Time To Repair (MTTR) of 30 hr. for failed nodes 12, the theoretical MTBF of the entire clique then becomes very high (a million time improvement). This is the fundamental benefit of redundancy. Thus, the parallel system is thought of as being composed of an array of cliques, rather than an array of individual processor nodes 12.

[0020] Figure 6 shows that a configuration can be held constant even though one or more processor nodes 12 has failed or are shut down. Data storage devices 14 have some well-known forms of redundancy and recovery (e.g., mirroring and other RAID algorithms). Configuring processor nodes 12 and data storage/devices 14 into cliques provide the redundant hardware, and vprocs 18 provide the redundant software. By using the traits of multiple threads 16 per processor node 12 and location independence (because a vproc 18 does not know which processor node 12 it is allocated to beforehand), it can be shown that the vproc 18 configuration can be held constant, even though one or more processor nodes 12 are unavailable. All of the data that is being processed is available on the shared data storage devices 14. The resultant configuration may operate with less processing power, but the parallel application still has the same number of vprocs 18 and retains full access to all data storage devices 14.

[0021] In summary, the present invention discloses a virtual processor method and apparatus for parallel computer systems that increases the level of parallelism to include multiple threads per node. If a processor node has a plurality of storage devices attached, each indi-

vidual thread can be allocated to a different data storage device. Similarly, if a processor node has multiple CPUs, each individual thread can utilize a different CPU in the node. Thus, a task could potentially occupy all available hardware in the system. The result is increased system utilization and availability. Other benefits derived are better control over the degree of parallelism, higher system availability without undue programming overhead in the application, and enhanced fault tolerance.

10

Claims

1. A computer system, including a plurality of processor nodes (12), each of the processor nodes (12) having a plurality of data storage devices (14) connected thereto; and an interconnect network (10), coupled to the processor nodes (12) for transmitting messages between processor nodes (12), each of the processor nodes (12) comprising a plurality of virtual processor means (18) for performing a plurality of threads (16) simultaneously, wherein each virtual processor means (18) performs one of the threads (16), the threads (16) being streams of instructions executed on behalf of a task, and each of the processor nodes (12) comprising operating system means for assigning the virtual processor means (18) to processor nodes (12) at system startup, **characterized in that** each of the threads (16) accesses data on a different data storage device (14) simultaneously with the other threads; and in that it further comprises means for creating and storing a translation table (34) identifying locations of virtual processor means (18) and processor nodes (12) on the interconnect network (10) at system startup.
2. A computer system according to claim 1, wherein each processor node(12) consists of one central processing unit (CPU).
3. A computer system according to claim 1, wherein each processor node (12) has a plurality of central processing units (CPUs).
4. A computer system according to any one of the preceding claims, comprising an interconnect network driver (28) for transmitting messages between virtual processor means (18), wherein each of the virtual processor means (18) includes one or more mailboxes, such that when a message is sent to a mailbox at a particular virtual processor means (18), the interconnect network driver (28) looks up a virtual processor means/processor node (12) translation in the translation table (34) and sends the message to a designated processor node and the designated processor node (12) routes it locally to an appropriate mailbox in the particular virtual proces-

- sor means (18).
5. A computer system according to any one of the preceding claims, comprising a clique formed by a clustering of processor nodes (12) connected to a common set of data storage devices (14), wherein the processor nodes (12) are cross-connected to the common set of data storage devices (14).
10. A method according to claim 8 or claim 9, including the steps of clustering processor nodes (12) into a clique connected to a common set of data storage devices (14), wherein the processor nodes (12) are cross-connected to the common set of data storage devices (14); and connecting each virtual processor means (18) to at least one virtual disk comprising one or more memory regions of the data storage devices (14) that are grouped together, such that the virtual disk can be accessed by a processor node (12) within a clique, thereby allowing the virtual processor means (18) to be started and executed on any processor node within a clique.
15. A computer system according to claim 5, wherein each virtual processor means (18) is connected to at least one virtual disk comprising one or more memory regions of the data storage devices (14) that are grouped together, wherein the virtual disk can be accessed by a processor node (12) within a clique, thereby allowing the virtual processor means (18) to be started and executed on any processor node (12) within a clique.
20. A computer system, according to any one of claims 4 to 6, wherein said operating system means are adapted to provide updates to the translation table when movement of the virtual processor means between processor nodes (12) occurs during (34) runtime.
25. A method of executing instructions in a parallel computer system comprised of a plurality of processor nodes (12), each of the processor nodes having a plurality of data storage devices (14) connected thereto, and an interconnect network (10), coupled to the processor nodes (12), for transmitting messages between processor nodes, comprising the steps of: utilizing a plurality of virtual processor means (18) to perform multiple threads (16) simultaneously, wherein each virtual processor means (18) comprises an operating environment in one of the processor nodes (12) for performing one of the threads (16), the threads (16) being streams of instructions executed on behalf of a task, and assigning each of the virtual processor means (18) to one of the processor nodes (12) at system startup, **characterized in that** each of the threads (16) accesses data on a different data storage device (14) simultaneously with the other threads (16); and further **characterized by** the steps of creating and storing a translation table (34) identifying locations of virtual processor means (18) and processor nodes (12) on the interconnect network (10) at system startup.
30. Computersystem mit einer Mehrzahl von Prozessorknoten (12), an deren jeden eine Mehrzahl von Datenspeichereinheiten (14) angeschlossen ist, mit einem Verbindungsnetzwerk (10) das mit den Prozessorknoten (12) gekoppelt ist, um Nachrichten zwischen diesen zu übertragen, wobei jeder Prozessorknoten (12) eine Mehrzahl virtueller Prozessoreinheiten (18) zur gleichzeitigen Bearbeitung einer Mehrzahl von Programmbausteingruppen (16) aufweist und wobei jeder virtuelle Prozessor (18) eine der Programmbausteingruppen (16) bearbeitet und die Programmbausteingruppen (16) durch Ströme von bezüglich einer Aufgabe auszuführenden Befehlen gebildet werden, wobei ferner jeder der Prozessorknoten (12) ein Betriebssystem zur Zuordnung der virtuellen Prozessoreinheiten (18) zu den Prozessorknoten (12) beim Systemanlauf enthält, **dadurch gekennzeichnet, dass** jede der Programmbausteingruppen (16) gleichzeitig mit anderen Programmbausteingruppen auf Daten in einer unterschiedlichen Datenspeichereinheit (14) zugreift, und dass sie Mittel zur Schaffung und Speicherung einer Zuordnungs-Tabelle (34) enthält, welche beim Systemanlauf Orte virtueller Prozessoreinheiten (18) und Prozessorknoten (12) in dem Verbindungsnetz (10) bezeichnet.
35. 2. Computersystem nach Anspruch 1, bei welchem jeder Prozessorknoten (12) aus einer Zentralprozessoreinheit (CPU) besteht.
40. 3. Computersystem nach Anspruch 1, bei welchem jeder Prozessorknoten (12) eine Mehrzahl zentraler Prozessoren (CPUs) aufweist.
45. 4. Computersystem nach einem der vorstehenden
- (34) and sends the message to a designated processor node (12) and the designated processor node (12) routes it locally to an appropriate mailbox in the particular virtual processor means (18).
5. (34) and sends the message to a designated processor node (12) and the designated processor node (12) routes it locally to an appropriate mailbox in the particular virtual processor means (18).
10. (34) and sends the message to a designated processor node (12) and the designated processor node (12) routes it locally to an appropriate mailbox in the particular virtual processor means (18).
15. (34) and sends the message to a designated processor node (12) and the designated processor node (12) routes it locally to an appropriate mailbox in the particular virtual processor means (18).
20. Patentansprüche
1. Computersystem mit einer Mehrzahl von Prozessorknoten (12), an deren jeden eine Mehrzahl von Datenspeichereinheiten (14) angeschlossen ist, mit einem Verbindungsnetzwerk (10) das mit den Prozessorknoten (12) gekoppelt ist, um Nachrichten zwischen diesen zu übertragen, wobei jeder Prozessorknoten (12) eine Mehrzahl virtueller Prozessoreinheiten (18) zur gleichzeitigen Bearbeitung einer Mehrzahl von Programmbausteingruppen (16) aufweist und wobei jeder virtuelle Prozessor (18) eine der Programmbausteingruppen (16) bearbeitet und die Programmbausteingruppen (16) durch Ströme von bezüglich einer Aufgabe auszuführenden Befehlen gebildet werden, wobei ferner jeder der Prozessorknoten (12) ein Betriebssystem zur Zuordnung der virtuellen Prozessoreinheiten (18) zu den Prozessorknoten (12) beim Systemanlauf enthält, **dadurch gekennzeichnet, dass** jede der Programmbausteingruppen (16) gleichzeitig mit anderen Programmbausteingruppen auf Daten in einer unterschiedlichen Datenspeichereinheit (14) zugreift, und dass sie Mittel zur Schaffung und Speicherung einer Zuordnungs-Tabelle (34) enthält, welche beim Systemanlauf Orte virtueller Prozessoreinheiten (18) und Prozessorknoten (12) in dem Verbindungsnetz (10) bezeichnet.
50. 2. Computersystem nach Anspruch 1, bei welchem jeder Prozessorknoten (12) aus einer Zentralprozessoreinheit (CPU) besteht.
55. 3. Computersystem nach Anspruch 1, bei welchem jeder Prozessorknoten (12) eine Mehrzahl zentraler Prozessoren (CPUs) aufweist.

- Ansprüche mit einem Verbindungsnetz-Treiber (28) zur Übermittlung von Daten zwischen virtuellen Prozessoreinheiten (18), deren jede eine oder mehrere Mailboxen enthält, so dass dann, wenn eine Nachricht an eine Mailbox einer speziellen virtuellen Prozessoreinheit geschickt wird, der Verbindungsnetz-Treiber (28) in der Zuordnungs-Tabelle (34) eine virtuelle Prozessoreinheit/Prozessorknoten (12) -Zuordnung aufsucht und die Nachricht an einen angegebenen Prozessorknoten schickt, der sie örtlich an eine richtige Mailbox in der virtuellen Prozessoreinheit (18) weiterleitet.
- 5
- 10
- 15
- 20
- 25
- 30
- 35
- 40
- 45
- 50
- 55
- und Zuordnung jeder der virtuellen Prozessoreinheiten (18) zu einem der Prozessorknoten (12) beim Systemanlauf, **dadurch gekennzeichnet**, dass jede der Programmabsteingruppen (16) auf Daten in einer anderen Datenspeichereinheit (14) gleichzeitig mit anderen Programmabsteingruppen (16) zugreift, und weiterhin **gekennzeichnet durch** die Schritte der Schaffung und Speicherung einer Zuordnungs-Tabelle (34), welche beim Systemstart Orte virtueller Prozessoreinheiten (18) und Prozessorknoten (12) in dem Verbindungsnetzwerk (10) identifiziert.
9. Verfahren nach Anspruch 8 mit dem Schritt der Übertragung von Nachrichten zwischen virtuellen Prozessoreinheiten (18), deren jede eine oder mehrere Mailboxen enthält, so dass beim Senden einer Nachricht an eine Mailbox einer bestimmten virtuellen Prozessoreinheit (18) ein Verbindungsnetz-Treiber nach einer virtuellen Prozessor/Prozessorknotenzuordnung in der Zuordnungstabelle (34) sucht und die Nachricht an einen angegebenen Prozessorknoten (12) schickt, welcher sie örtlich an eine passende Mailbox in der speziellen virtuellen Prozessoreinheit (18) schickt.
10. Verfahren nach Anspruch 8 oder 9 mit den Schritten: Häufung von Prozessorknoten (12) zu einer Clique, die mit einem gemeinsamen Satz von Datenspeichereinheiten (14) verbunden ist, wobei die Prozessorknoten (12) mit dem gemeinsamen Satz von Datenspeichereinheiten (12) kreuzverbunden sind, und Anschließen jeder virtuellen Prozessoreinheit (18) an mindestens einen virtuellen Plattspeicher, der einen oder mehrere Speicherbereiche der zu Gruppen zusammen gefassten Datenspeichereinheiten (14) umfaßt, so dass auf den virtuellen Plattspeicher von einem Prozessorknoten (12) innerhalb der Clique zugegriffen werden kann und die virtuelle Prozessoreinheit (18) an irgendeinem Prozessorknoten (12) innerhalb der Clique gestartet und durchgeführt werden kann.
7. Computersystem nach einem der Ansprüche 4 - 6, bei welchem die Betriebssysteme geeignet sind, Aktualisierungen der Zuordnungs-Tabelle durchzuführen, wenn während (34) der Laufzeit virtuelle Prozessoreinheiten zwischen den Prozessorknoten (12) wechseln.
8. Verfahren zur Ausführung von Befehlen in einem Parallelcomputersystem aus einer Mehrzahl von Prozessorknoten (12), deren jeder mit einer Mehrzahl von Datenspeichereinheiten (14) verbunden ist, und mit einem mit den Prozessorknoten (12) gekoppelten Verbindungsnetzwerk (10) zur Übertragung von Nachrichten zwischen den Prozessorknoten mit den Schritten:
- Benutzung einer Mehrzahl virtueller Prozessoreinheiten (18) zur gleichzeitigen Bearbeitung mehrerer Programmabsteingruppen (16), wobei jede virtuelle Prozessoreinheit (18) eine Betriebsumgebung in einem der Prozessorknoten (12) aufweist zur Bearbeitung einer der Programmabsteingruppen (16), welche durch Ströme von Befehlen gebildet werden, die bezüglich einer Aufgabe auszuführen sind,

Revendications

1. Système informatique comportant une pluralité de noeuds processeurs (12), chacun des noeuds processeurs (12) ayant une pluralité de dispositifs (14) de stockage de données reliés à celui-ci ; et un réseau d'interconnexion (10), relié aux noeuds processeurs (12) pour transmettre des messages entre les noeuds processeurs (12), chacun des noeuds processeurs (12) comprenant une pluralité de moyens processeurs virtuels (18) pour exécuter simultanément une pluralité de processus (16), dans

- lequel chaque moyen processeur virtuel (18) exécute l'un des processus (16), les processus (16) étant des suites d'instructions exécutées pour le compte d'une tâche, et chacun des noeuds processeurs (12) comprenant des moyens formant système d'exploitation pour affecter les moyens processeurs virtuels (18) à des noeuds processeurs (12) au démarrage du système, **caractérisé en ce que** chacun des processus (16) accède à des données se trouvant sur un dispositif de stockage de données (14) différent en même temps que les autres processus ; et en ce qu'il comprend en outre des moyens pour créer et stocker une table de conversion (34) identifiant les emplacements des moyens processeurs virtuels (18) et des noeuds processeurs (12) sur le réseau d'interconnexion (10) au démarrage du système.
2. Système informatique suivant la revendication 1, dans lequel chaque noeud processeur (12) consiste en une unité centrale de traitement (CPU).
3. Système informatique suivant la revendication 1, dans lequel chaque noeud processeur (12) comporte une pluralité d'unités centrales de traitement (CPU).
4. Système informatique suivant l'une quelconque des revendications précédentes, comprenant un pilote (28) de réseau d'interconnexion destiné à transmettre des messages entre des moyens processeurs virtuels (18), dans lequel chacun des moyens processeurs virtuels (18) comporte une ou plusieurs boîtes aux lettres, de telle façon que lorsqu'un message est envoyé à une boîte aux lettres d'un moyen processeur virtuel (18) particulier, le pilote (28) du réseau d'interconnexion consulte une conversion entre un moyen processeur virtuel et un noeud processeur (12) dans la table de conversion (34) et envoie le message à un noeud processeur désigné, et le noeud processeur (12) désigné l'achemine localement à une boîte aux lettres particulière dans le moyen processeur virtuel (18) particulier.
5. Système informatique suivant l'une quelconque des revendications précédentes, comprenant une clique formée d'un regroupement de noeuds processeurs (12) connectés à un ensemble commun de dispositifs de stockage de données (14), les noeuds processeurs (12) étant interconnectés à l'ensemble commun de dispositifs de stockage de données (14).
6. Système informatique suivant la revendication 5, dans lequel chaque moyen processeur virtuel (18) est connecté à au moins un disque virtuel comprenant une ou plusieurs régions de mémoire dans les dispositifs de stockage de données (14) qui sont re-
- 5 groupés les uns avec les autres, le disque virtuel pouvant faire l'objet d'un accès par un noeud processeur (12) d'une clique, ce qui permet aux moyens processeurs virtuels (18) d'être démarrés et exécutés sur un noeud processeur quelconque d'une clique.
- 10 7. Système informatique suivant l'une quelconque des revendications 4 à 6, dans lequel lesdits moyens formant système d'exploitation sont aptes à fournir des mises à jour de la table de conversion lorsqu'un mouvement des moyens processeurs virtuels entre des noeuds processeurs (12) se produit pendant (34) le temps d'exécution.
- 15 8. Procédé d'exécution d'instructions dans un système informatique parallèle constitué d'une pluralité de noeuds processeurs (12), chacun des noeuds processeurs ayant une pluralité de dispositifs de stockage de données (14) qui lui sont connectés et un réseau (10) d'interconnexion, relié aux noeuds processeurs (12) pour transmettre des messages entre des noeuds processeurs (12) comprenant les étapes consistant à : utiliser une pluralité de moyens processeurs virtuels (18) pour exécuter simultanément des processus multiples (16), chaque moyen processeur virtuel (18) comprenant un environnement d'exploitation dans l'un des noeuds processeurs (12) pour exécuter l'un des processus (16), les processus (16) étant des suites d'instructions exécutées pour le compte d'une tâche, et à affecter chacun des moyens processeurs virtuels (18) à l'un des noeuds processeurs (12) au démarrage du système, **caractérisé en ce que** chacun des processus (16) accède à des données se trouvant sur un dispositif de stockage de données (14) différent en même temps que les autres processus (16) ; et **caractérisé par** les étapes consistant à créer et à stocker une table de conversion (34) identifiant des emplacements de moyens processeurs virtuels (18) et de noeuds processeurs (12) sur le réseau d'interconnexion (10) au démarrage du système.
- 20 30 35 40 45 50 55 9. Procédé suivant la revendication 8, comportant l'étape consistant à transmettre des messages entre des moyens processeurs virtuels (18), chacun des moyens processeurs virtuels (18) comportant une ou plusieurs boîtes aux lettres, de telle sorte que lorsqu'un message est envoyé à une boîte aux lettres dans un moyen processeur virtuel (18) particulier, un pilote du réseau d'interconnexion consulte une conversion entre un processeur virtuel et un noeud processeur dans la table de conversion (34) et envoie le message à un noeud processeur (12) désigné et que le noeud processeur (12) désigné l'achemine localement à une boîte aux lettres appropriée dans le moyen processeur virtuel (18)

particulier.

10. Procédé suivant la revendication 8 ou 9, comprenant les étapes consistant à regrouper des noeuds processeurs (12) sous la forme d'une clique connectée à un ensemble commun de dispositifs de stockage de données (14), dans lequel les noeuds processeurs (12) sont interconnectés à l'ensemble commun de dispositifs de stockage de données (14) ; et à connecter chaque moyen processeur virtuel (18) à au moins un disque virtuel comprenant une ou plusieurs régions de mémoire des dispositifs de stockage de données (14) qui sont regroupés les uns avec les autres, de telle sorte que le disque virtuel puisse faire l'objet d'un accès par un noeud processeur (12) d'une clique, ce qui permet aux moyens processeurs virtuels (18) d'être démarrés et exécutés sur un noeud processeur (12) quelconque d'une clique.

5

10

15

20

25

30

35

40

45

50

55

FIG. 1

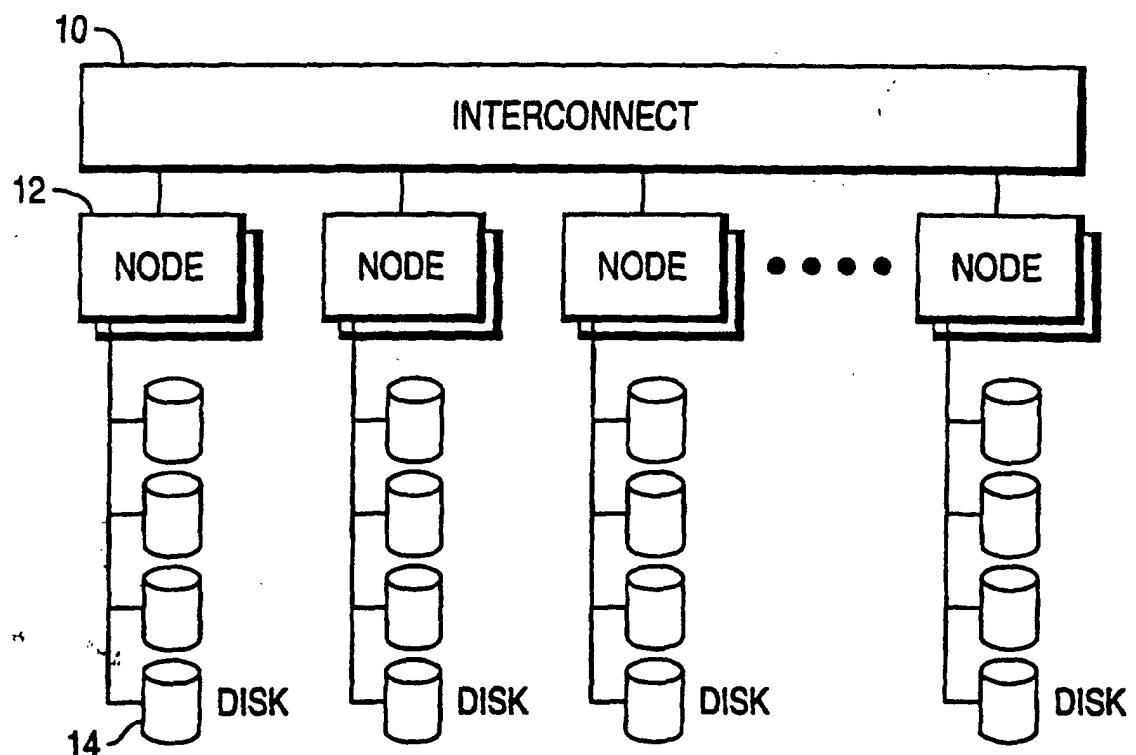


FIG. 2

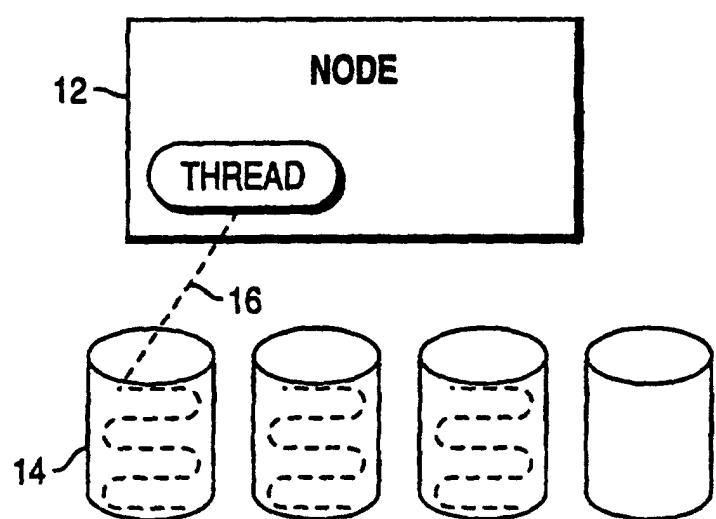


FIG. 3

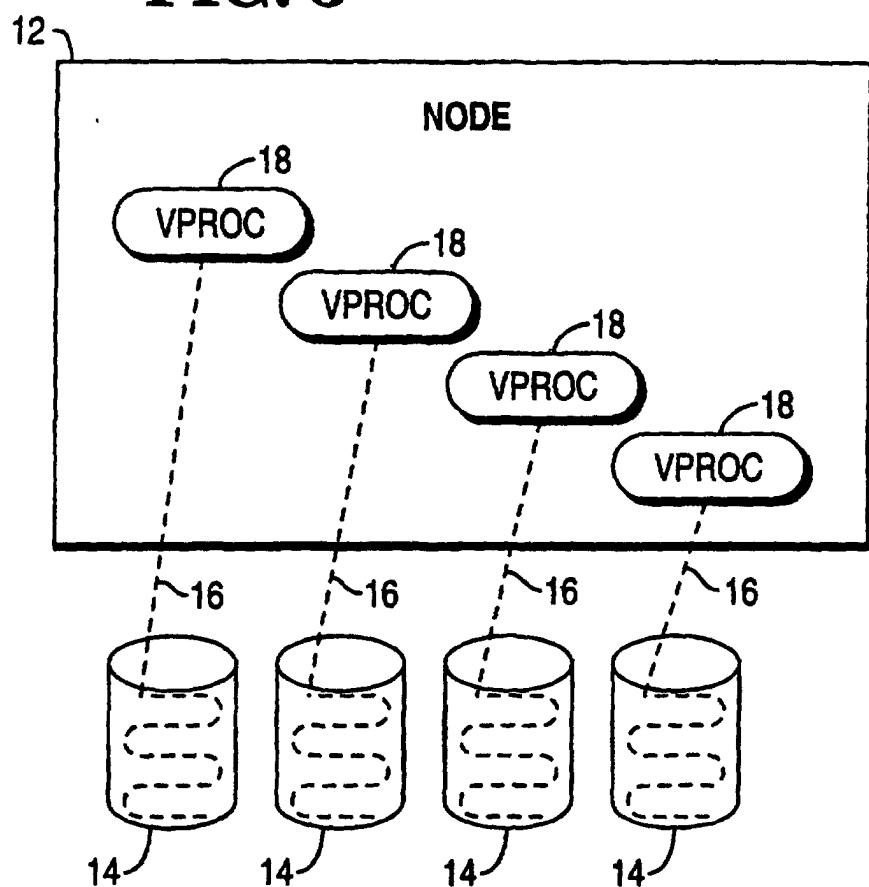


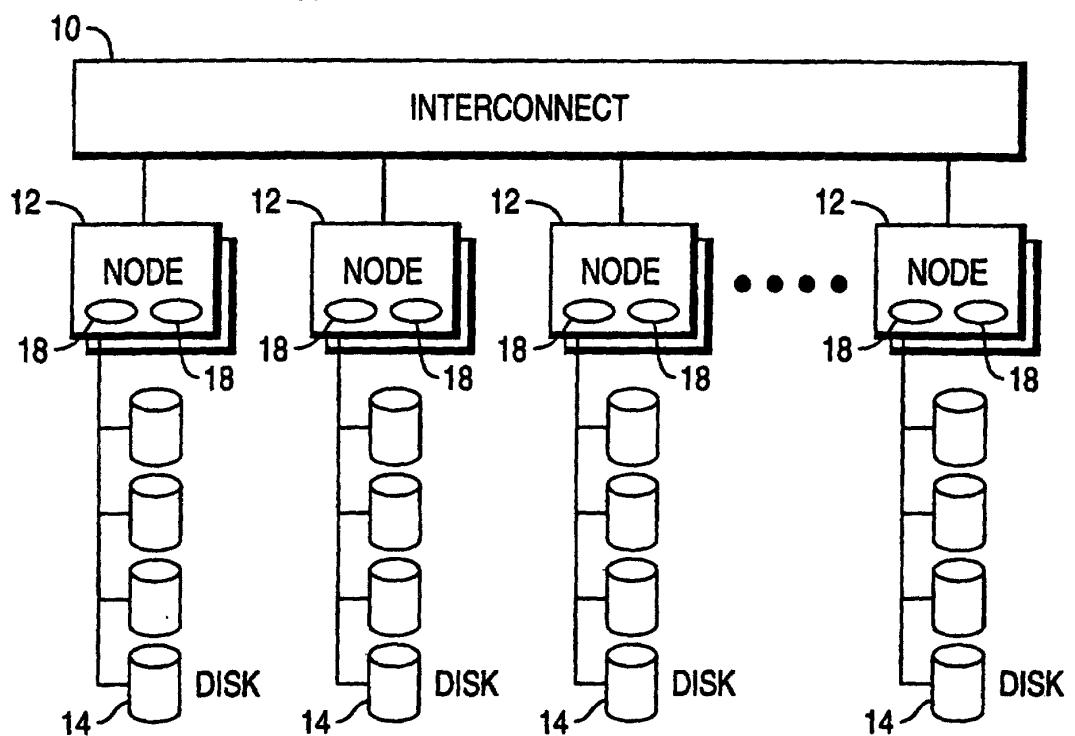
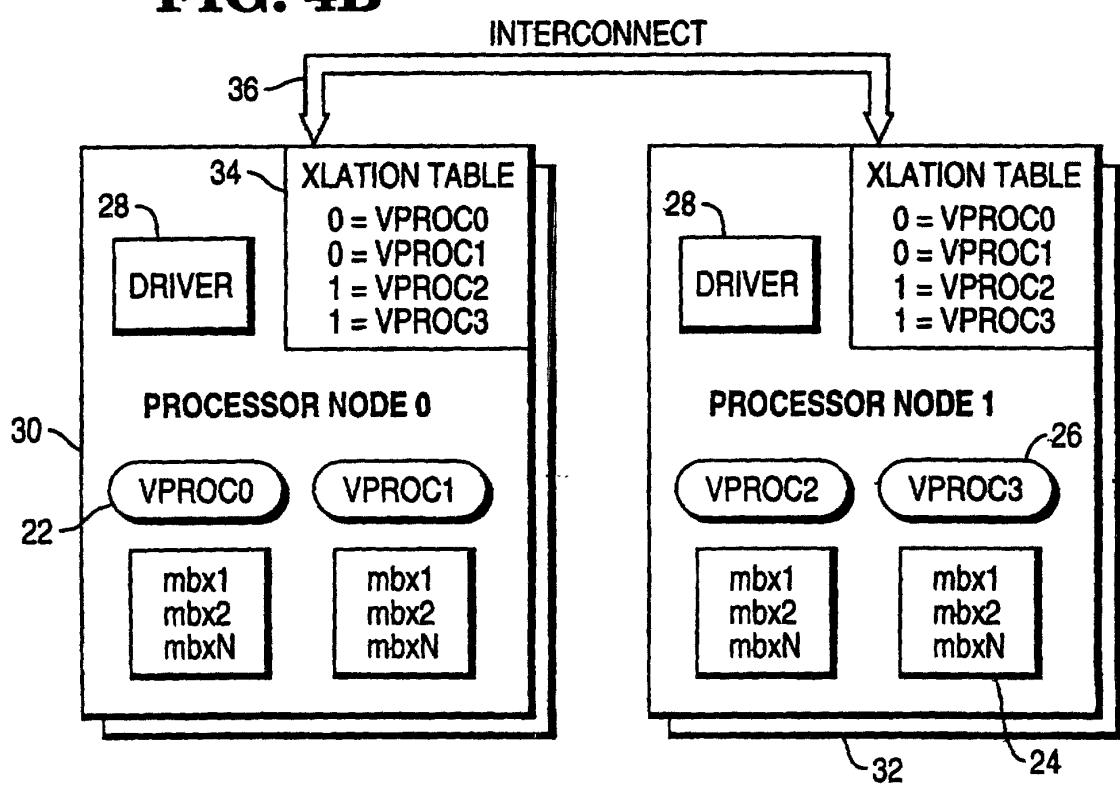
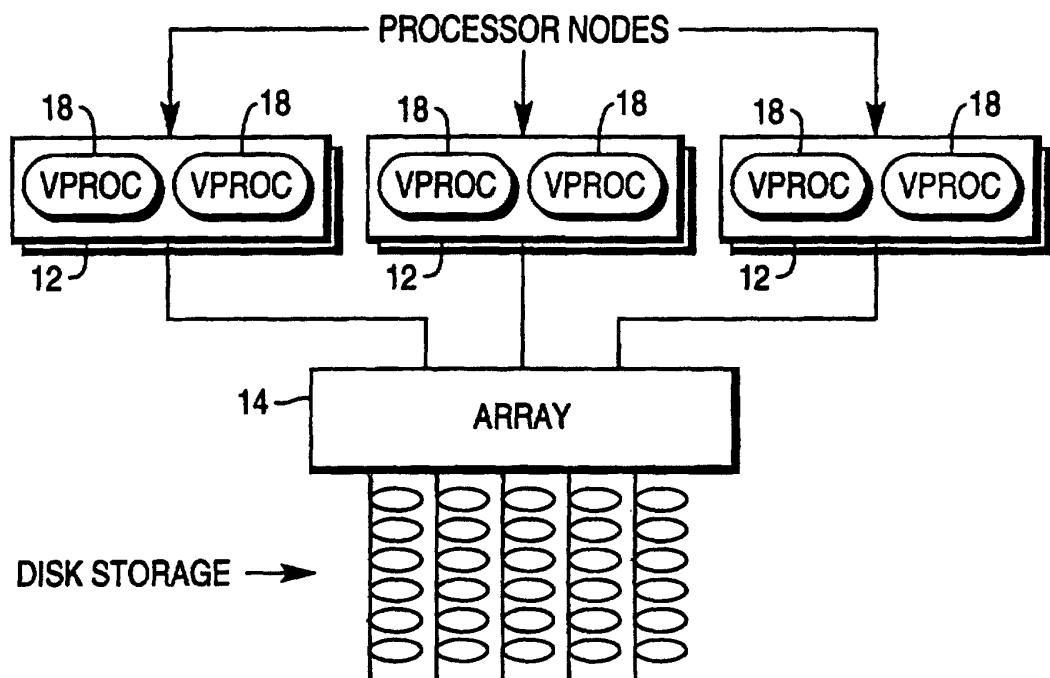
FIG. 4A**FIG. 4B**

FIG. 5**FIG. 6**